

## 2003/09/24

(43)公開日 平成11年(1999)12月24日

(51)Int.Cl. <sup>6</sup>	識別記号	F I	
G 0 6 F 15/16		G 0 6 F 15/16	4 2 0 J
	3 1 0	11/14	3 1 0 Z
	3 1 0	11/16	3 1 0 Z

審査請求 未請求 請求項の数 4 O.L (全 9 頁)

(21)出願番号	特願平10-161808	(71)出願人	000005108 株式会社日立製作所 東京都千代田区神田駿河台四丁目6番地
(22)出願日	平成10年(1998)6月10日	(72)発明者	本堂 友理 神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所システム開発研究所内
		(72)発明者	長須賀 弘文 神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所システム開発研究所内
		(72)発明者	秋葉 千江子 神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所ソフトウェア開発本部内
		(74)代理人	弁理士 小川 勝男

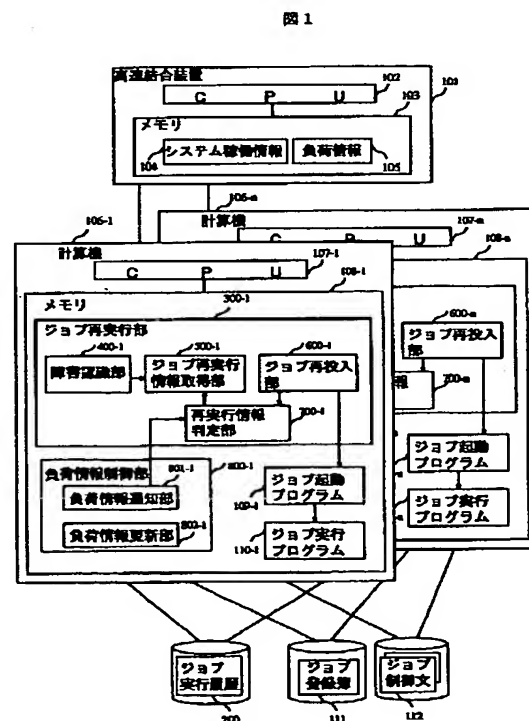
最終頁に続く

(54) 【発明の名称】 ジョブ再実行方法

(57) 【要約】

【課題】 複数の計算機から構成される計算機システムにおける計算機障害時にジョブを再実行させる場合の、負荷分散機能を用いた効率的なジョブ再実行技術を提供する。

【解決手段】障害の発生した計算機で実行中であつたジョブを再実行させる際に、再実行を行う候補となる計算機の負荷情報と、被再実行ジョブの中断されるまでの実行状況の情報から、再実行させる計算機を判定することで効率的な実行と計算機障害の影響の早期解決を実現する。



**【特許請求の範囲】**

【請求項 1】複数の計算機から構成される計算機システムに、上記計算機システム上で動作している全てのジョブが使用する計算機資源と稼働状況に関する情報を収集し、上記システム内の全ての計算機ノードが上記収集した情報を参照することを可能とするステップを有する計算機システムのジョブの再実行方法において、第一の上記計算機ノードに障害が発生した場合に、上記計算機システム内の当該第一の計算機を除いた計算機から成る第二の計算機群の各々は、上記収集した情報から上記第一の計算機上で実行していたジョブの特性を解析し、判断基準に基づいて自計算機ノードで再実行させるジョブを上記第一の計算機上で実行していたジョブの中から求め、求めたジョブを自計算機ノード上で再実行させるステップを有することを特徴とするジョブ再実行方法。

【請求項 2】請求項 1 記載のジョブ再実行方法において、上記判定基準として CPU 使用率を用い、自計算機の CPU 負荷が他計算機に対して低く、自計算機の CPU 負荷から得る CPU 使用率の閾値を越えない CPU 使用率であることが上記収集した情報から得られるジョブを求めて再実行を行うことを特徴とするジョブ再実行方法。

【請求項 3】請求項 1 記載のジョブ再実行方法において、上記判定基準として上記中断したジョブが CPU を多く使用するジョブであるか又は入出力処理を多く行うジョブであるかを用い、自計算機の CPU 負荷が低い場合は CPU を多く使用するジョブを求め、自計算機の入出力装置の使用率が低い場合は入出力処理を多く行うジョブを求めて再実行を行うことを特徴とするジョブ再実行方法。

【請求項 4】請求項 1 記載のジョブ再実行方法において、上記判定基準として上記中断したジョブの中断するまでの実行時間を用い、自計算機の予め指定された閾値を越えない実行時間を持つジョブを求めて再実行を行うことを特徴とするジョブ再実行方法。

**【発明の詳細な説明】****【0001】**

【発明の属する技術分野】本発明は、複数計算機からなり、各々の計算機の負荷情報が相互に取得可能な手段を有する計算機システムでのジョブ再実行方法に関し、特にジョブ実行中の計算機の障害発生時におけるジョブ再実行方法に関するものである。

**【0002】**

【従来の技術】大規模システムでは、その社会的用途から高性能及び高信頼性が求められており、その要求から生まれた機能の一つに、システム障害時におけるジョブの自動再実行機能がある。複数計算機から成る計算機システムにおいて、ある計算機に障害が発生し実行継続が不可能となったジョブの再実行方式の一例として特開平

7-175766 号公報の「疎結合多重システムのジョブ再実行制御方式」がある。

【0003】これは、ジョブを投入する場合は予め再実行させる計算機を定義しておき、ある計算機に障害が発生した場合は、その計算機の回復を待たずに、上記の予め定義しておいた別計算機にスケジューリングすることで、速やかに再実行が行えるような手段を提供するものである。

**【0004】**

【発明が解決しようとする課題】大規模システムでは、複数の計算機をノードとして接続することで構成されるクラスタ型計算機システム上で業務実施するのが主流になりつつある。このようなクラスタ型計算機システム上で効率よく業務処理を遂行するためには、各ノードの計算機資源を有効に活用する必要がある。そのためには、特定の計算機ノードに業務処理が集中することがないように、各計算機ノードに処理を分散させている。

【0005】上記従来技術は、障害により再実行する必要の生じたジョブをスケジュールする計算機は、予めジョブ制御言語や特定の情報格納領域に記述されていることが前提となっている。従って、上記のようなクラスタ型計算機システムで上記従来技術を実現した場合は、障害が発生した場合は特定の計算機に処理が集中し、効率よい業務処理遂行が妨げられる可能性がある。

【0006】また、上記従来技術では、障害が発生した場合には、特定の計算機がマスタ計算機としてジョブの再スケジュールリングを行うことになっている。従って、そのマスタ計算機自体に障害が発生した場合は回復作業に支障をきたす可能性がある。上記クラスタ型計算機システムは、従来のマスタスレーブ構成をとる疎結合マルチプロセサシステムとは異なり、特定の計算機をマスタとして固定することなく自律的に処理を行い耐障害性を高めることを特徴とするシステムである。従って固定的な特定のマスタ計算機が存在を前提とする上記従来技術ではこのクラスタ型計算機システムの利点を損なう恐れがあった。

【0007】本発明は複数計算機から成りシステムを構成する各計算機の負荷情報を取得する手段を有する計算機システムにおけるジョブの再実行方法に関するものであり、以下の二つの目的を持つ。

【0008】本発明の第一の目的は、この負荷情報取得手段を利用し、再実行すべきジョブの特性をふまえた再スケジュールリングすることで、計算機障害の影響を抑えた効率的なクラスタ型計算機システムの運用を実現することにある。

【0009】本発明の第二の目的は、複数の計算機ノードがジョブの回復手段を有することで、特定の計算機ノードの障害によるジョブ再実行不能状態を回避し、システム全体の耐障害性を向上させることにある。

**【0010】**

【課題を解決するための手段】効率よく業務処理を遂行するためには、特定の計算機ノードに業務処理が集中することがないように、各計算機ノードに処理を分散させなければならない。そのために、クラスタ型計算機システムでは各々の計算機の負荷を収集し、その情報をシステム全体の処理状況として各計算機ノードで実行される処理が取得できる負荷情報取得手段が具備されているものが多い。本発明ではこの負荷情報取得手段を利用し、再実行すべきジョブをその特性に見合った計算機に再スケジューリングすることで、計算機障害の影響を抑えたクラスタ型計算機システムの運用を実現する。

【0011】ジョブを実行する複数の計算機とそれらの負荷情報等を格納するために上記複数の計算機各々から書き込み・参照可能な領域を保持する1つの高速結合装置、または共有拡張記憶装置、または特定計算機の記憶装置、またはその他の記憶装置から成り、上記複数の計算機が自計算機の負荷情報を収集しシステム全体の負荷情報を格納する領域に対する更新を行う負荷情報更新処理機能と、上記格納された負荷情報を要求に応じて通知する負荷情報通知処理機能からなる負荷情報制御手段を有し、ジョブの実行内容を定義したジョブ制御文を格納した記憶装置と、ジョブの実行順序を登録するジョブ登録簿を格納した記憶装置と、各計算機で実行しているジョブの実行履歴を格納した記憶装置がシステム内の計算機各々から共通して参照及び更新可能である計算機システムで、上記第一の目的を達成するために、上記各計算機からアクセス可能な領域におかれた計算機システムの稼働情報からシステム内の計算機に発生した障害を検知するステップと、上記検知ステップにおいて障害が発生した場合は、発生した計算機とそこで実行されていたジョブの情報を取得し、そのジョブを再実行する準備をするステップと、上記準備をするステップの通知により、障害の発生した計算機での実行履歴から再実行するジョブの情報を取得するステップと、ジョブを再実行する際のシステムの各計算機ノードの負荷情報を取得するステップと、上記ジョブの情報と計算機ノードの負荷情報から、上記ジョブを再実行する計算機を決定するステップと、上記決定により自計算機が上記ジョブを再実行するに適した計算機であると判定された場合は、そのジョブを自計算機に再投入し実行を行うステップを設ける。ここで、上記障害の発生した計算機での実行履歴から得るジョブの情報は、ジョブの中断までの実行時間、実行時間に占めるCPU使用時間、前記CPU使用時間と前記実行時間の比率、使用した入出力装置台数等が挙げられる。

【0012】上記第二の目的を達成するために、上記計算機システム内の各々の計算機が、上記第一の目的を達成するための上記複数のステップを具備し、ある計算機に障害が発生した場合は、他の計算機上の上記複数のステップが実行されるものとする。

【0013】上記各計算機からアクセス可能な領域におかれた計算機システムの稼働情報からシステム内の計算機に発生した障害を検知するステップと、上記検知ステップにおいて障害が発生した場合は、発生した計算機とそこで実行されていたジョブの情報を取得し、そのジョブを再実行する準備をするステップにより、ジョブの再実行を行う計算機が自律的に中断したジョブの再投入を行うことができる。

【0014】上記準備をするステップの通知により、障害の発生した計算機での実行履歴から再実行するジョブの情報を取得するステップと、ジョブを再実行する際のシステムの各計算機ノードの負荷情報を取得するステップにより、自計算機で再実行するためのジョブの情報と、計算機システム内での自計算機の負荷情報を把握するための情報を取得することができる。

【0015】上記ジョブの情報とノードの負荷情報から、上記ジョブを再実行する計算機を決定するステップと、上記決定により自計算機が上記ジョブを再実行するに適した計算機であると判定された場合は、そのジョブを自計算機に再投入し実行を行うステップにより、再実行すべきジョブの特性に合わせた計算機での再実行が可能となる。

【0016】

【発明の実施の形態】本発明の実施形態を図を用いて詳細に説明する。

【0017】まず第一の実施形態について図1～図4を用いて説明する。

【0018】本実施形態では、一つ以上の複数の計算機により構成されるシステムを対象とする。複数の計算機は相互に通信を行う手段として、各計算機から更新・参照可能である領域を保持する高速結合装置により各計算機が接続されているものとする。また、ジョブの実行内容を定義したジョブ制御文を格納した記憶装置と、ジョブの実行順序を登録するジョブ登録簿を格納した記憶装置と、各計算機で実行しているジョブの実行履歴を格納した記憶装置が各計算機から共通して更新・参照可能であるものとする。

【0019】図1は、本実施形態の基本的な構成を表している。CPU102およびメモリ103から構成される高速結合装置101により接続された複数の計算機106は各々がCPU107およびメモリ108から構成される。

【0020】システム内の各計算機の負荷情報を管理する負荷情報制御部800は計算機106のメモリ108上に置かれ、負荷情報更新部802と負荷情報通知部801から構成される。負荷情報更新部802は一定時間毎に自計算機の負荷情報を収集し、その内容を高速結合装置101のメモリ103上の負荷情報105の領域に書き込む処理を行う。負荷情報通知部801は、本発明のジョブ再実行部300やその他のプログラムの要求に

より、負荷情報 105 の情報を通知したり、負荷情報 105 の内容から最も負荷の低い計算機を判定し結果を通知したりする処理を行う。

【0021】各計算機 106 は自計算機が稼働しているという情報をシステム稼働情報 104 に置き、他計算機に障害が発生したか否かという情報を取得できるようにしている。

【0022】各々の計算機 106 のメモリ 108 上には障害により中断されたジョブを再実行するジョブ再実行部 300 が設置されている。ジョブ再実行部 300 は、システム稼働情報 104 から他計算機の障害を認識する障害認識部 400 と、再実行すべきジョブの情報をジョブ実行履歴 200 とジョブ制御文 112 から取得して再実行の準備を行うジョブ再実行情報取得部 500 と、取得したジョブの情報と負荷情報通知部 801 から取得した負荷情報 105 からジョブの特性を判定し、再実行するのに適した計算機を選択する再実行情報判定部 700 と、再実行するのに適した計算機が自計算機であった場合は自計算機に対してジョブの再投入を行うジョブ再投入部 600 から構成される。ここでは、中断したジョブの情報として、そのジョブの中断するまでの平均 CPU 使用率を用いて説明を行う。

【0023】各々の計算機 106 のメモリ 108 上には、ジョブ再投入部 600 からの通知でジョブの起動を行うジョブ起動プログラム 109 と、ジョブの実行を行うジョブ実行プログラム 110 が設置されている。

【0024】ジョブ実行履歴 200 の構造を図 2 に示す。ジョブ実行履歴 200 は計算機 106 のジョブ実行単位である空間に対応した複数のレコードから構成される。その空間で実行されていたジョブ名、実行を開始した日付、開始した時刻、そのジョブの障害が発生した時刻までの CPU 使用時間等が含まれている。

【0025】障害により実行が中断されたジョブの再実行を行うジョブ再実行部 300 の処理の流れを図 3 により説明する。

【0026】他計算機に障害が発生したか否かの情報を取得し、障害が発生した場合は他計算機で実行が中断されたジョブの回復処理を開始する（ステップ 400）。

【0027】障害の発生した計算機の障害が発生した時刻等の情報を含むシステム稼働情報 104 と、障害の発生した計算機で実行されていたジョブの情報を含むジョブ実行履歴 200 と、ジョブの処理内容を定義したジョブ制御文 112 からジョブの情報を取得する（ステップ 500）。

【0028】再実行すべきジョブが存在するかを判定する（ステップ 301）。上記判定が偽であった場合、処理を終了する。上記判定が真であった場合、取得した情報から再実行情報判定処理 700 により、再実行するに適した計算機を求める（ステップ 302）。

【0029】上記判定から得られた再実行に適した計算

機が自計算機であるか判定を行う（ステップ 303）。上記判定が偽であった場合、ステップ 304 から処理を行う。上記判定が真であった場合、自計算機のジョブ起動プログラム 109 に対し、ジョブの起動要求を行う（ステップ 600）。

【0030】まだ再実行すべきジョブがあるかを判定する（ステップ 304）。上記判定が真であった場合、ステップ 302 から次の処理を行う。上記判定が偽であった場合、処理を終了する。

【0031】再実行情報判定部の処理の流れを図 4 により説明する。

【0032】ジョブ再実行プログラムから再実行すべきジョブの情報を取得する。システム稼働情報 104 に含まれる障害が発生した時刻の情報と、ジョブ実行履歴 200 に含まれるジョブが実行開始した時刻の情報から、中断するまでの経過時間が取得できる。また、ジョブ実行履歴 200 からそのジョブの中断するまでの CPU 使用時間が取得できる。ここから対象となるジョブの平均した CPU 使用率が求められる（ステップ 701）。

【0033】次に負荷情報通知部 801 から、再実行するジョブが実行可能な計算機の現時点での負荷情報を取得する（ステップ 701）。

【0034】ステップ 701 で求めたジョブが消費する CPU 使用率を割り当て可能な計算機を判定する（ステップ 703）。

【0035】ステップ 703 により求めたジョブを再実行するのに適した計算機の判定結果を戻す（ステップ 704）。

【0036】以上、第一の実施形態を具体的に説明したが、前記実施形態において再実行するジョブから得られる情報をジョブの中断するまでの平均 CPU 使用率としたが、これを中断するまでのジョブの実行時間としてもよい。その場合は、長大な実行時間がかかるジョブを特定計算機に分担させるといったシステムの形態が実現できる。

【0037】また、再実行するジョブから得られる情報を、ジョブ制御文 112 から得られる入出力装置の台数としてもよい。その場合は、入出力装置の台数に比較的余裕のある計算機に振り分けることが可能となる。

【0038】以上、本発明の実施形態において各計算機が接続される装置を高速結合装置として説明したが、その装置を共有拡張記憶装置としてもよい。その場合の基本的な構成を図 5 に示す。システムを構成する全ての計算機 106-n と接続されている共有拡張記憶装置 113 上にシステム稼働情報 104 と負荷情報 105 が設置されるものとする。

【0039】また、上記の領域を保持する装置をディスク装置としてもよい。その場合の基本的な構成を図 6 に示す。システムを構成する全ての計算機 106-n と接続されているディスク装置 114 上にシステム稼働情報

1 0 4 と負荷情報 1 0 5 が設置されるものとする。

【0 0 4 0】

【発明の効果】本発明によれば、システム内の計算機の障害によって再実行しなければならないジョブが短時間に大量に発生しても、効率的に処理できるようにジョブの特性にあわせてシステム内に分散させ再実行することができる。

【0 0 4 1】さらに、障害のため実行が中断したジョブを再実行するため分散させる機能を特定計算機に偏らせることなく、その時の計算機システムの運用状況に合わせて分散することが可能となり、高性能化、高信頼化でできる。

【図面の簡単な説明】

【図 1】本発明のジョブ再実行方法の第一の実施例を示した計算機システムの構成図である。

【図 2】本発明の入力情報となるジョブ実行履歴の説明図である。

【図 3】ジョブ再実行部の処理を説明したフローチャートである。

【図 4】再実行情報判定部の処理を説明したフローチャートである。

【図 5】本発明のジョブ再実行方法を実現した計算機システムの別の構成図である。

【図 6】本発明のジョブ再実行方法を実現した計算機システムの別の構成図である。

【符号の説明】

1 0 1…高速結合装置、 1 0 2…高速結合装置の CPU、1 0 3…高速結合装置のメモリ、1 0 4…システム稼働情報管理テーブル、1 0 5…負荷情報、  
1 0 6 - n…計算機、1 0 7 - n…計算機の CPU、1 0 8 - n…計算機のメモリ、3 0 0 - n…再実行部、  
4 0 0 - n…障害認識部、5 0 0 - n…ジョブ再実行情報取得部、6 0 0 - n…ジョブ再投入部、7 0 0 - n…再実行情報判定部、8 0 0 - n…負荷情報制御部、1 0 9 - n…ジョブ起動プログラム、  
1 1 0 - n…ジョブ実行プログラム、2 0 0 - n…ジョブ実行履歴、1 1 1 - n…ジョブ登録簿、  
1 1 2 - n…ジョブ制御文。

【図 2】

図 2

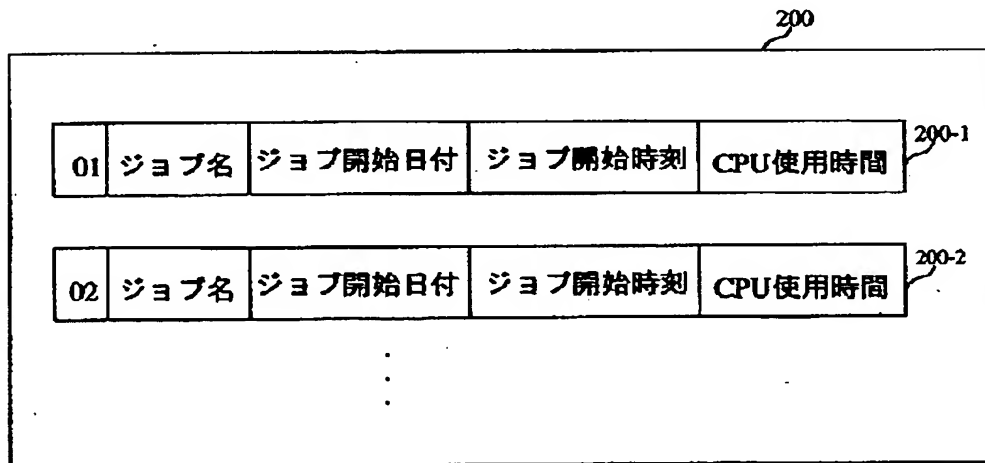
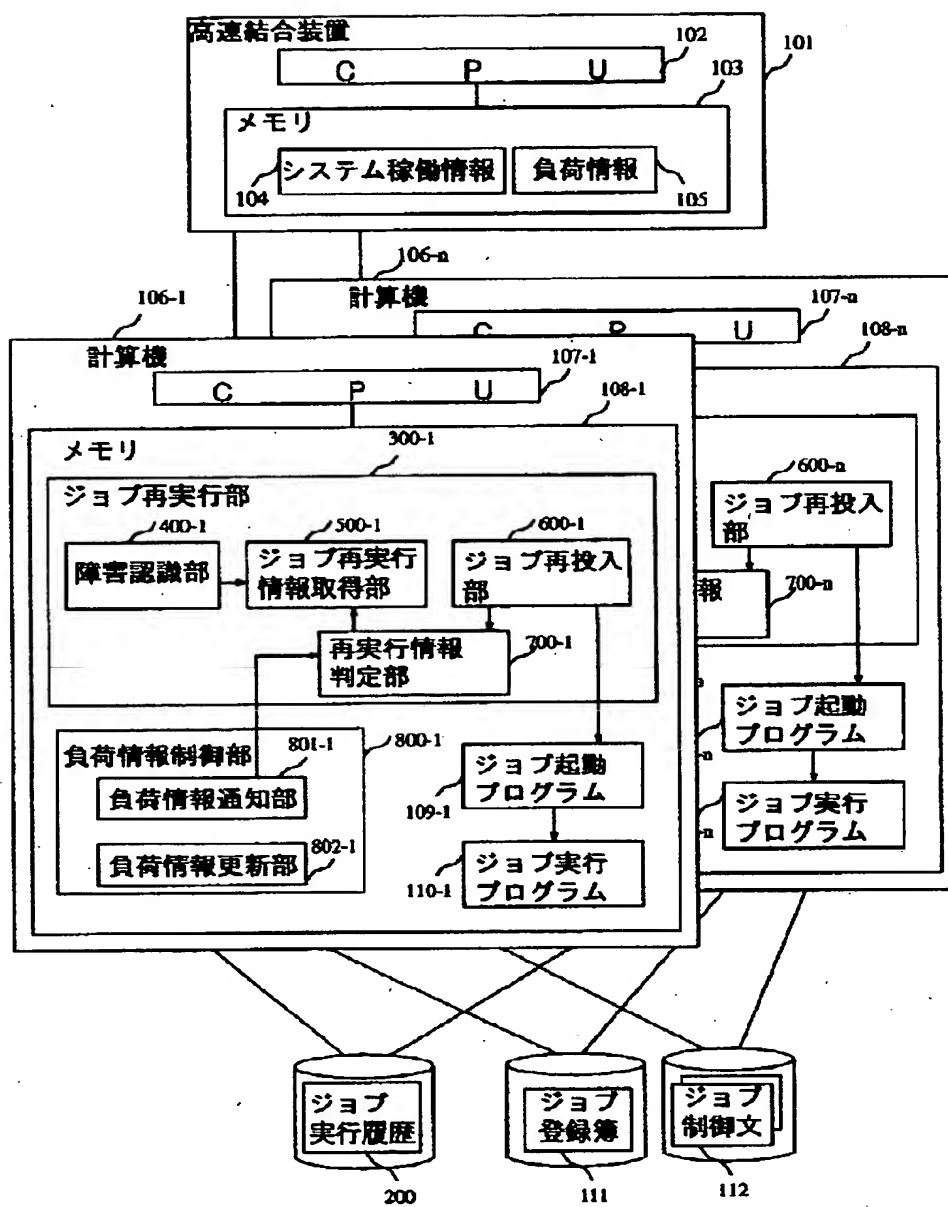
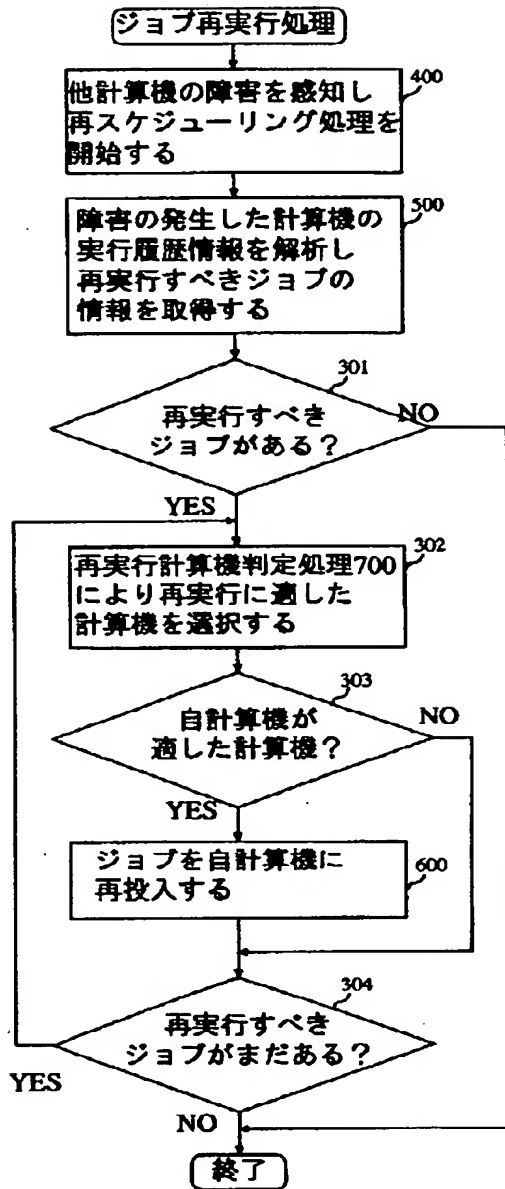


Figure 1 shows a schematic diagram of a single-layered structure. It consists of a rectangular block with a wavy top surface. A horizontal line is drawn across the middle of the block. Below this line, there is a label '1'.



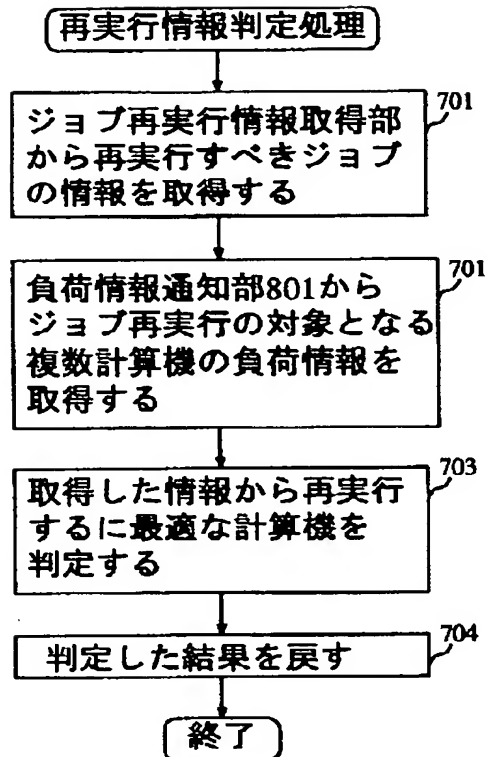
【図 3】

図 3



【図 4】

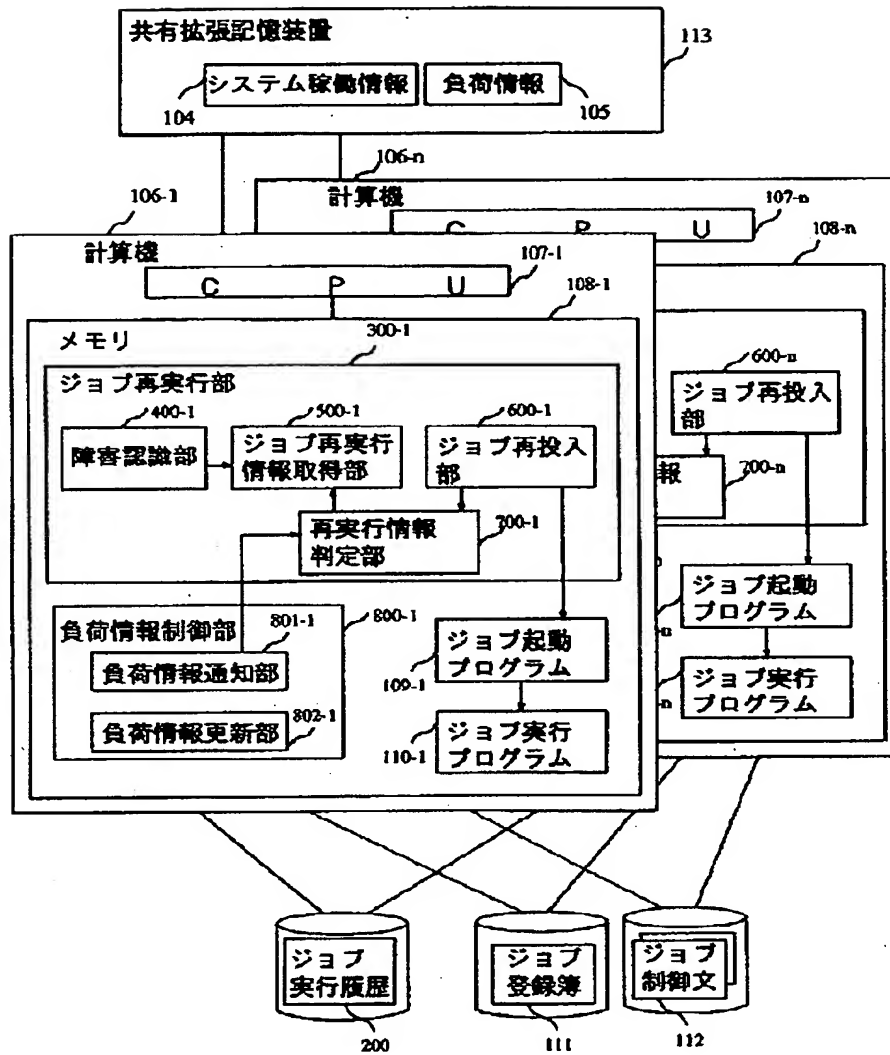
図 4





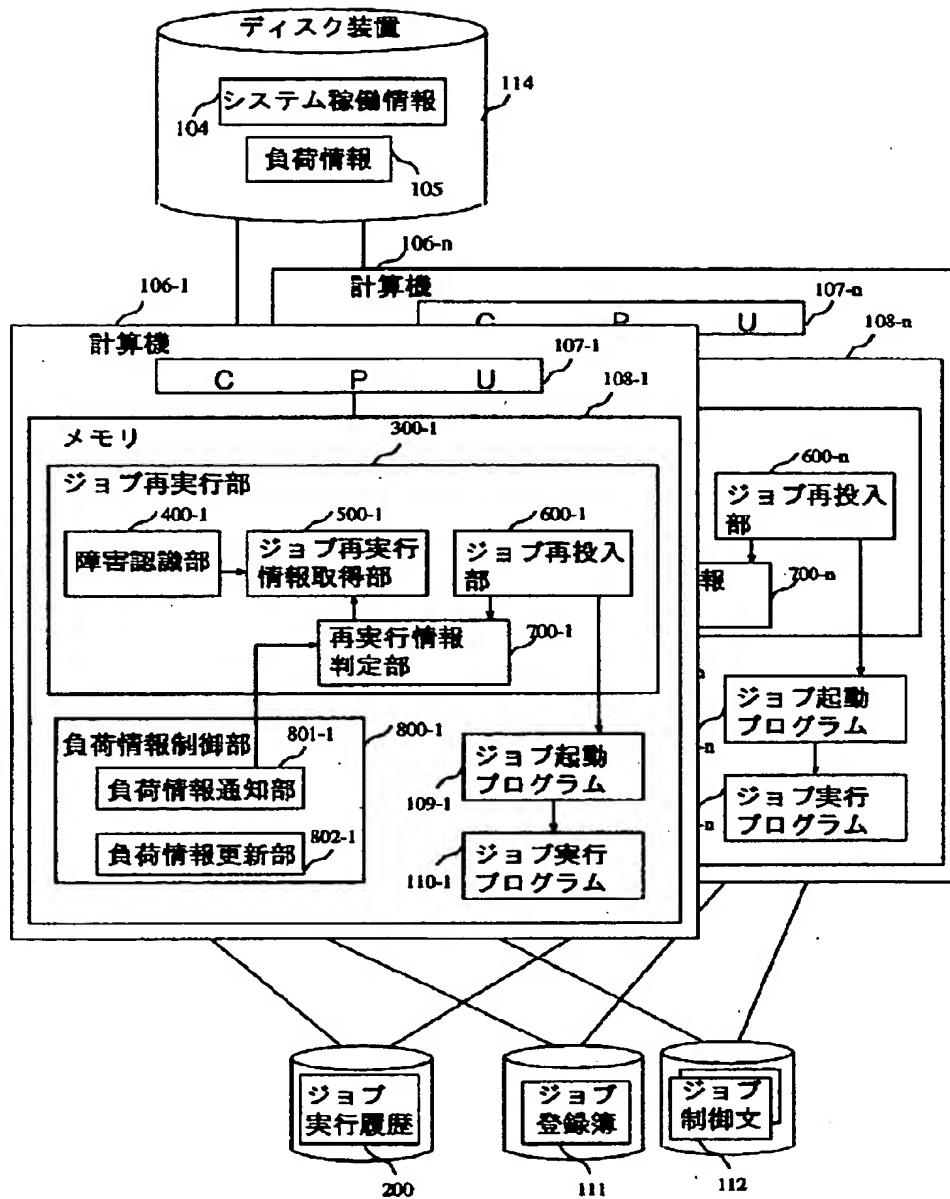
【図5】

図5



【図6】

図6



フロントページの続き

(72)発明者 岩倉 義之  
神奈川県横浜市戸塚区戸塚町5030番地 株  
式会社日立製作所ソフトウェア開発本部内